

AM JA999118

日 本 国 特 許 庁

PATENT OFFICE  
JAPANESE GOVERNMENT



別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日  
Date of Application:

1999年 7月 7日

出 願 番 号  
Application Number:

平成11年特許願第192618号

出 願 人  
Applicant (s):

インターナショナル・ビジネス・マシーンズ・コーポレイション

CERTIFIED COPY OF  
PRIORITY DOCUMENT

1999年11月26日

特許庁長官  
Commissioner,  
Patent Office

近 藤 隆 彦



出証番号 出証特平11-3083336

【書類名】 特許願

【整理番号】 JA999118

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 15/00

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 東京基礎研究所内

【氏名】 那須川 哲哉

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 東京基礎研究所内

【氏名】 長野 徹

【特許出願人】

【識別番号】 390009531

【住所又は居所】 アメリカ合衆国 1 0 5 0 4、ニューヨーク州アーモンク（番地なし）

【氏名又は名称】 インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】 100086243

【弁理士】

【氏名又は名称】 坂口 博

【連絡先】 0 4 6 2 - 7 3 - 3 3 1 8、3 3 2 5、3 4 5 5

【選任した代理人】

【識別番号】 100091568

【弁理士】

【氏名又は名称】 市位 嘉宏

【手数料の表示】

【予納台帳番号】 024154

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9304391

【包括委任状番号】 9304392

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 データ分析システム

【特許請求の範囲】

【請求項 1】

データから有効な知識を取り出す、データ分析システムであって、

- (1) 文書データを含むデータから概念を抽出する、概念抽出手段と、
- (2) 前記抽出された概念から特徴的な概念を抽出する、特徴的概念抽出手段と

を含む、データ分析システム。

【請求項 2】

前記概念抽出手段(1)が、文書データを含むデータからカテゴリ別の概念を抽出する手段であり、前記特徴的概念抽出手段(2)が、前記カテゴリ別の概念において、同一カテゴリに属する概念のうち、対応する別のカテゴリに属する概念の中で占める割合が既定値を超えている概念を抽出する手段である、請求項 1 記載のシステム。

【請求項 3】

前記概念抽出手段(1)が、

データ中の前記文書データを形態素解析する手段と、

前記形態素解析の結果に基づき、前記文書データの文節を生成する手段と、

前記文節内のキーワードの概念を抽出する手段であって、前記文節に対してカテゴリ辞書を適用して、文節内のキーワードに対してカテゴリを付加する手段と、前記文節を有する文を構文木生成規則に従い構文木を生成する、構文解析手段と

前記文節内のカテゴリを付加されたキーワードに関し、同一文中でのキーワードの係り受け関係を抽出する、係り受け抽出手段と、

前記カテゴリ別の概念を抽出する手段であって、前記キーワード間の係り受け関係に基づき、係り受けの関係にある各キーワードのカテゴリの組合せを抽出する手段と、

を含む、請求項 2 記載のシステム。

【請求項 4】

前記特徴的な概念を検出する手段（2）が、  
ユーザの命令を受け取る入力手段と、前記ユーザの命令を解析する、命令解析手段と、前記解析された命令に従い、前記カテゴリ別の概念を提示し、同一カテゴリに属する概念のうち、対応する別のカテゴリに属する概念の中で占める割合が既定値を超えている概念を他の概念と異なる属性で表示する手段とを有する、請求項 2 記載のシステム。

【請求項 5】

前記特徴的な概念を検出する手段（2）が、さらに、  
抽出された概念の相対頻度を計算する手段と、  
抽出された概念からキーワードを検索する手段と、  
カテゴリ別キーワードの頻度を計算する手段と、  
得られた前記相対頻度、検索結果、キーワードの頻度を表示する手段と、  
を含む、請求項 4 記載のシステム。

【請求項 6】

データから有効な知識を取り出す、概念抽出方法であって、  
データ中の文書データを形態素解析を行う段階と、  
前記形態素解析の結果に基づき、前記文書データの文節を生成する段階と、  
前記文節に対してカテゴリ辞書を適用して、文節内のキーワードに対してカテゴリを付加する段階と、  
前記文節を有する文を構文木生成規則に従い構文木を生成する段階と、  
前記文節内のカテゴリを付加されたキーワードに関し、同一文中でのキーワードの係り受け関係を抽出する段階と、  
前記キーワード間の係り受け関係に基づき、係り受けの関係にある各キーワードのカテゴリの組合せを抽出する段階と、  
を有する、概念抽出方法。

【請求項 7】

データから有効な知識を取り出すプログラムを含む媒体であって、該プログラムが、コンピュータに、

データ中の文書データを形態素解析を行う機能と、  
 前記形態素解析の結果に基づき、前記文書データの文節を生成する機能と、  
 前記文節に対してカテゴリ辞書を適用して、文節内のキーワードに対してカテゴリを付加する機能と、  
 前記文節を有する文を構文木生成規則に従い構文木を生成する機能と、  
 前記文節内のカテゴリを付加されたキーワードに関し、同一文中でのキーワードの係り受け関係を抽出する機能と、  
 前記キーワード間の係り受け関係に基づき、係り受けの関係にある各キーワードのカテゴリの組合せを抽出する機能と、  
 を実現させる、プログラムを含む媒体。

【発明の詳細な説明】

【0001】

【産業上の利用分野】

本発明は、大量のデータから知識を獲得する技術に関し、特にキーワードに付加した同一カテゴリの概念比較により、大量の非定形文からなるデータから注目に値する有効な知識を獲得する方法およびシステムに関する発明である。

【0002】

【従来の技術】

大量の文書の内容を様々な形で視覚化し分析に役立てようという試みが次第に増えてきている。しかし、従来の方法では、名詞句を中心としたキーワードを抽出し、キーワードの分布を表示（例えば、キーワード間の相関関係を計算し、その結果に基づいて相関の強いキーワードをクラスタ化して表示）する程度の処理しか行っていない。その程度の処理の結果を用いた場合には、ユーザーが様々な観点からデータを絞り込みながら、主観的に注目に値すると思える部分を発見するまで手探りで分析するしかなく、分析過程におけるユーザーの負担が大きい。また、多種多様なキーワードを混合して、まとまりの無いままの状態で扱うため、有効な結果にたどり着くのは難しい。

つまり、テキストをデータマイニングの手法（クラスタリングや相関ルールの分析など）で解析してみようという試みは、データマイニングが注目され始めた

頃から存在するものの、従来の手法では、テキストから抽出した分析単位が、単なる文字列としてのキーワードに過ぎなかった為、有効な結果を得られずに終わってしまうケースが多い。

【 0 0 0 3 】

【発明が解決しようとする課題】

従って、本発明が解決しようとする課題は、大量データから有効な知識を獲得する方法およびシステムを提供することである。

また別の課題は、大量データから注目に値する知識を自動的に見つけ出す方法およびシステムを提供することである。

また別の課題は、大量のデータから有効な知識を獲得するにあたり、ユーザビリティに優れた分析方法およびシステムを提供することである。

【 0 0 0 4 】

【課題を解決するための手段】

上記課題を解決するために、本発明は文書を含む大量のデータから特異な特徴を有する概念を抽出することにより、有効な知識を獲得する方法およびシステムを提供する。

本発明は、概念抽出装置と特徴的概念抽出装置と有する。前記概念抽出は文書データを含むデータからカテゴリ別の概念を抽出する。前記特徴的概念抽出装置は前記抽出した概念の中から特徴的な概念を抽出する装置であって、前記カテゴリ別の概念において、同一カテゴリに属する概念のうち、対応する別のカテゴリに属する概念の中で占める割合が既定値を超えている概念を抽出する。

より、詳細には、概念抽出装置では、語彙辞書や文法知識を利用した形態素解析やカテゴリ辞書を用いた曖昧性解消手法を用いて、非定型テキストから、カテゴリ別の概念を抽出する。特徴的概念抽出装置では、カテゴリとカテゴリの各組み合わせにおいて、同じカテゴリの概念のうち、対応する別のカテゴリの概念の占める割合の比や差が既定値（しきい値）を超えているものを見つけて出す。さらにカテゴリとカテゴリの各組み合わせを表形式で表示し、該表中において注目すべき項目を他と異なる属性で表示したり、リストアップして提示するなどの処理を行う。

## 【0005】

## 【発明の実施の形態】

図1は、本発明のデータ分析システム概要を示すブロック図である。なお実施例としてコンピュータ製品に関する電話の問い合わせデータに基づき、これを解析し、特徴ある概念を抽出するデータ分析システムについて説明する。図1のブロック110は問い合わせデータ150を入力としてラベル付きデータ160を出力するデータ変換部である。問合せデータ150に含まれる非定型のデータと、同じく含まれる定型のデータを同一に保持する形式のデータを作成する。ブロック120は、入力された文に対してカテゴリ辞書170を用いてカテゴリを付加し、カテゴリを付加されたキーワードのうち、同じ文中で係り受け関係のあるものを、より具体的な意味を表現する概念（ラベル・概念付きデータ180）として抽出する、概念抽出部である。ブロック130は、ラベル・概念付きデータ180を入力として、特徴のある概念を検索・抽出する検索・特徴検出部130である。このブロックは検索・特徴抽出を効率良く行えるようにするため、データ全体もしくは部分集合における頻度分布などの統計情報を作成・保持する。同じカテゴリに属する概念は同じような振る舞いをする（同じような出現傾向を持ち、同じような概念と共起する）傾向があると想定し、各概念の振る舞いを全体、もしくは部分集合の値と比較することで、効率良く、注目すべき情報を検索・検出する。さらにこの検索・特徴検出部130は、統計情報を視覚的に表示する機能を有し、特徴ある概念の分布差異の提示を行う。

## 【0006】

上記各ブロック110、120、130について以下に詳細に説明する。

## 〔データ変換部110〕

まずブロック110に入力される問い合わせデータの例は以下のようなものである。

----- 問い合わせデータ例 -----

1999/01/01

0000001

タイトル：ノートで日本語が使えない



マシンタイプ：製品A

問題種別：総合案内

CALL種別：案内

回答・対応種別：窓口対応

解決期間：1日

通話時間：21分

ご質問：ノートパソコンで日本語が使えないので、OSを再インストールしたのですが、それ以降、モデムとイーサネットカードが使えなくなっていました。

----- 問い合わせデータ例 --終わり-----

このように、問合せデータは日付や定型項目からの選択肢、また連続値や離散値、自由に記入できるタイトルや文章など、さまざまな形式の項目から構成されたデータからなる。

【0007】

ブロック110は上記のような問い合わせデータ150を入力とし、下記のようなラベル付きデータに変換する。

----- ラベル付きデータ -----

ID199901010000001

TITPで日本語が使えない

KWM1MT：製品A

KWQ3TC：総合案内

KWQ4TD：案内

KWQ2PT：窓口対応

KWP3SD：1日

KWP4CM：21分

CTQ：ノートパソコンで日本語が使えないので、OSを再インストールしたのですが、それ以降、モデムとイーサネットカードが使えなくなっていました。

----- ラベル付きデータ --終わり---

このように、非定形データを含む問い合わせデータを上記のようなラベル付きデ

ータに変換することにより、さまざまな型のデータを同じ形式に変換する。上の例ではTIはタイトル、MTはマシンタイプ、TCは問題種別、TDはCALL種別、PTは回答・対応種別、SDは解決期間、CMは通話時間、CTは元の問合せ内容、KW+2byteは項目の種類を示す。図2にデータ変換部のフローチャートを示す。ステップ210で問い合わせデータ150を読み込み、ステップ220でデータの終了であるかを判断し、そうでなければステップ230で形式を変換する。データが終了した場合はステップ240で変換を終了する。

【0008】

〔概念抽出部120〕

図3にブロック120の概念抽出部を説明する。ブロック310は、ラベル付きデータ160内の文に対して、形態素解析を行う、形態素解析装置である。次にブロック320は、ブロック310で形態素に分解された文について形態素を文（あるいは特定の文脈）に出現する順に文節生成規則に従って文節を決定する、文節生成装置である。形態素が付属語である場合や明らかに文節が切れると判断されたところで文節を区切りこれを文末になるまで行う。

【0009】

次にブロック330は、ブロック320で生成された文節に対してカテゴリ辞書340を適用し文節内のキーワードに対してカテゴリを付加する辞書適用装置である。キーワードにカテゴリを付加することで単なる文字列ではなく意味を持った概念として扱う。（例えば、「ワシントン」という文字列を単なるキーワードとして扱えば、人名も地名も区別できず、有効な分析が行えないが、[人名][地名]といったカテゴリを付加することで、意味を持つことが出来る）体言（名詞類の語）は（元表現 品詞 概念 カテゴリ）のデータで表現されるカテゴリ辞書を参照してカテゴリを付加する。述語に関しては体言と同様にカテゴリ辞書を用いると共に、付属語の情報から[問題][要望][疑問]といったカテゴリ付けを行う。例えば、「壊れた」という動詞は（壊れる [動詞] 故障 [問題]）というデータがカテゴリ辞書にあれば、[問題]というカテゴリに属する「故障」という概念として抽出されるが、「…できない」「…したい」という表現は、カテゴリ辞書を参照せずに[問題][要望]として解釈することができる。上記カテゴリ辞書17

0 の構造の例を以下に示す。

----- カテゴリ辞書の例 -----

ノートパソコン

固有名詞

ノートパソコン

N1

O S

固有名詞

オペレーティングシステム

N2

----- カテゴリ辞書の例 - 終わり -

上記辞書においてカテゴリはN1（ハードウェア）、N2（ソフトウェア）というように対応付けられている。

【0 0 1 0】

ブロック 3 5 0 はブロック 3 3 0 でキーワードに対してカテゴリが付加された文節を有する文を簡単な構文木生成規則で構文木を生成する構文木解析装置である。

【0 0 1 1】

ブロック 3 6 0 は文節内のカテゴリを付加されたキーワードのうち同じ文中で係り受け関係のあるものを、より具体的な意味を表現する概念として抽出する係り受け抽出装置である。このブロック 3 6 0 では構文解析装置 3 5 0 の構文解析結果により得られるキーワード間の係り受け関係に基づき、係り受けの関係にある各キーワードのカテゴリの組合せを概念（ラベル・概念付きデータ 3 7 0）として抽出する。ラベル・概念付きデータ 3 7 0 の例を以下に示す。

----- ラベル・概念付きデータの例 -----

ID199901010000001

TI ノートで日本語が使えない

KWM1MT : 製品 A

KWQ3TC : 総合案内

KWQ4TD : 案内

KWQ2PT : 窓口対応

KWP3SD : 1 日

KWP4CM : 2 1 分

CTQ：ノートパソコンで日本語が使えないので、OSを再インストールしたのですが、それ以降、モデムとイーサネットカードが使えなくなっていました。

KWN1ノートパソコン

KWN0日本語

KWV2使えない

KWW6ノートパソコン…使えない

KWN2OS

KWV6再インストールする

KWWDOS…再インストールする

KWN1モデム

KWN1イーサネットカード

KWV2使えない

KWW6モデム…使えない

KWW6イーサネットカード…使えない

----- ラベル・概念付きデータの例 --終わり--

上記のようにラベル・概念付きデータ 3 7 0 は、データ変換部で得られたラベル付きデータ 1 6 0 に、概念抽出部 1 2 0 で抽出されたデータを追加した形を取り、ラベル付きデータ 1 6 0 と同じく同一形式のデータとなる。

【0 0 1 2】

図 1 2 に本発明の概念抽出部の流れを実際の文に基づき説明する。まずステップ 1 2 1 0 において、入力文”MODEMとイーサネットカードが使えない。”が入力された場合、ステップ 1 2 2 0 において、形態素解析装置 3 1 0 により、文を単語に区切り、品詞番号が付加される。例えば、

104…固有名詞、 81…格助詞「と」、 75…格助詞「が」、 10…動詞語幹、 44…???, 51…打ち消しの助動詞「ない」、100…句読点、などのような対応表を用いて形態素解析を行う。これにより入力文”MODEMとイーサネットカードが使えない。”は、以下のように変換される。

[MODEM, 104] [と, 81] [イーサネット, 104] [カード, 104] [が, 75] [使, 10] [え, 44] [ない, 51] [。 , 100]

【0013】

次にステップ1230、およびステップ1240で文節生成が行われる。ステップ1230の文節生成1では、形態素解析された文の単語列を文節にまとめる操作を行う。あらかじめ”{81, 75, 100, …}で文節を区切る”というルールを用意しておき、このルールを文頭から適用して、文頭から文節ごとに区切っていく操作を行う。上記入力文の場合、先頭から3文節あるが、各文節の先頭の単語がそれぞれ、名詞・名詞・動詞であることから、それぞれ順に、体言句・体言句・用言句であると判断される。その結果、入力文は以下のように変換される。

{[MODEM, 104] [と, 81]}  
 {[イーサネット, 104] [カード, 104] [が, 75]}  
 {[使, 10] [え, 44] [ない, 51] [。 ,100]}

【0014】

次にステップ1240の文節生成2では、文節生成1で区切られた文節に対して、すべての文節を自立語と付属語の組にする働きを行う。体言句は、名詞が複数含まれる場合は先頭から名詞を結合するようにする。例えば、{[イーサネット, 104] [カード, 104] → [イーサネットカード, 104]}のようにする。その後自立語の品詞コードを、一般名詞句を示すN1に書き換える。用言句は、付属語列（[え, 44] [ない, 51] [。 ,100]）を解析し、否定の情報を示す[ない, 51]を取り出し、動詞の語幹[使, 10]は[え, 44]と結合して、終止形「使える」にする。品詞コードを一般動詞を示すV1にし、否定の情報を付加して-V1とする。その結果、入力文は以下のように変換される。

{[MODEM, N1] [と, 81]}  
 {[イーサネットカード, N1] [が, 75]}  
 {[使える, -V1] [。 , 100]}

次にステップ1250で、カテゴリ辞書を用いて、自立語と付属語の組に分解された文節に対してカテゴリ付を行う。ここで適用される辞書は、以下の3つである。

(MODEM    N1    モデム    NA)  
 (イーサネットカード    N1    イーサネットカード    NA)

(使える -V1 使えない VC)

なおNAはハードウェアを意味し、VCは問題を意味する。その結果、入力文は以下のように変換される。

{[モデム, NA] [と, 81]}

{[イーサネットカード, NA] [が, 75]}

{[使えない, VC] [。 , 100]}

【0 0 1 5】

次にステップ 1 2 6 0 で、カテゴリ付された文節に文節からなる文に基づき構文木を生成する。この時の係り受けルールの形式は（係り受け元文節の自立語，係り受け元文節の付属語，係り受け先文節の自立語，係り受け先文節の付属語）となっている。このルールを文頭の文節1 {[モデム, NA] [と, 81]} から適用する。文節n番目に対してn+1番目から最終文節Nまで適用する(n = 1 ~ N-1)。係り受けルール中に (NA, 81, VC, \*) というルールがあるので、 {[モデム, NA] [と, 81]} と {[使えない, VC] [。 , 100]} に係り受けの関係があると判断される。なおルール中の\*はすべての品詞またはカテゴリーにマッチするという意味である。これを(n = 1 ~ N-1)で行い、係り受け情報を含む文節を有効グラフとして表し、(係り受け元, 係り受け先, 自立語, カテゴリー, 付属語の品詞番号) という形式に変換する。その結果、入力文は以下のように変換される。

(1, 3, "モデム", NA, 81)

(2, 3, "イーサネットカード", NA, 75)

(3, NULL, "使えない", VC, 100)

【0 0 1 6】

最後に、ステップ 1 2 7 0 において、構文解析された文を入力として、係り受け抽出抽出ルールに従って係り受けを抽出する。抽出ルールは任意の長さのカテゴリの列からなる。例えば(カテゴリ 1, カテゴリ 2, ..., カテゴリn) という形式である。文節番号1からNまでの係り受けを見て、係り受け抽出ルールに (NA, VC) というルール(n = 2)があるので、「モデム…使えない」と「イーサネットカード…使えない」という2つの係り受けが抽出される。

結局、元の文書である「モデムとイーサネットカードが使えない」という文章が

ら

「モデム …使えない」 [ハードウェア…問題]

「イーサネットカード …使えない」 [ハードウェア…問題]

という概念情報が取り出されたことになる。このようにして取り出された概念情報は、ラベル・概念付きデータベース 180 に登録される。

【0017】

図4に、本発明の概念抽出部 120 の処理のフローチャートを示す。

ステップ420でラベル付きデータ 160 内の文章 T を形態素  $W_0 \sim W_m$  に分割する。ここで形態素 W は、文字列  $w$  と品詞  $p$  で表される。すなわち  $W = \{w, p\}$  である。(以上が形態解析装置の処理である)

【0018】

次にステップ430で全単語が文節に変換されたかどうかを判断し、そうであれば処理はステップ440へ移り、そうでなければステップ432で単語  $W_n$  が付属語または句読点かどうかを判断し、その結果が No であればステップ434で、文節  $P_i$  に単語  $W_n$  を追加する。ここで文節 P は、1つ以上の連続した単語 W の集合で  $P = \{W\} = \{\{w, p\}\}$  である。そして処理はステップ430に戻る。ステップ432の判断の結果が Yes であればステップ436で文節  $P_i$  に単語  $W_n$  を追加した後、その次の文節を用意する。処理はその後ステップ440に移る。ステップ440で、全文節に対する処理が必要かどうかを判断し、必要であれば処理はステップ450へ移り、費用でなければステップ442において、文節を  $P = \{\{w, p\}\}$  の形式から  $P' = \{\{w_1, p_1\} \{w_2, p_2\}\}$  に変換する。

(ここで  $w_1$  は自立語  $w_2$  は付属語である) 例えば  $P = \{[国際, 名詞] [情勢, 名詞] [は, 助詞]\}$  であれば、名詞句はひとつにまとめ、 $P' = \{[国際情勢, 名詞] [は, 助詞]\}$  とする。そして処理はステップ440へ戻る。(以上が文節生成装置 320 の処理である)

【0019】

ステップ450では、自立語に対して辞書引きが終了したかどうかを判定する。その判定結果が No であれば、ステップ452で、 $((w_1, p_1) == (w_a, p_a))$  のエントリに対し  $(w_1 = w_b, p_1 = p_b)$  とする。同時に用言の処理を行う。この時

用いられる辞書 454 は [wa pa wb pb] のエントリの集合からなり、それぞれ [元表現, 品詞, 概念, カテゴリ] を表す。ここでの概念というエントリは、置き換え表現を意味する。元表現が「PC」であれば置き換え表現は「パソコン」となる。また全エントリはハッシュ構造で格納されているので高速にアクセスできる。例えばエントリ集合は、[マシン, 名詞, 機械, ハードウェア]、[壊れる, 動詞, 壊れる, 問題] 等である。この辞書を用いて、名詞句  $P' = \{ [マシン, 名詞] [が, 助詞] \}$  であれば、 $P'' = \{ [機械, ハードウェア] [が, 助詞] \}$  となる。用言の句  $P' = \{ [壊れる, 動詞] [ない, 助詞] \}$  であれば、動詞→問題となるはずだが、「ない」があるので、 $P'' = \{ [壊れない, 好評] [NULL, NULL] \}$  とする。また辞書に該当するエントリがなければ何もしない。ステップ 450 の判定結果が Yes であれば処理はステップ 460 へ進む。(以上が辞書適用装置 330 の処理である)

#### 【0020】

ステップ 460 で、構文木が完成したかどうか判定される。構文木が完成した場合、処理はステップ 470 へ進む。構文木が完成していない場合、処理はステップ 462 で、一般的に行われる構文解析を行い、結果として、 $P_n$  と  $P_k$  がリンクされる。(以上が構文解析装置 350 の処理である。)

#### 【0021】

ステップ 470 で、係り受けの抽出が終了したかどうかを判断する。もし終了していなければ処理はステップ 472 へ進み、ルールに基づき、 $P_n$  に対してリンクされた 2 項関係を抽出し、ラベル・概念付きデータベース 180 に登録する。この時、係り受けルール 474 を参照する。係り受けルール 474 のルールは [px py] のエントリの集合からなり、それぞれ [係り受け元カテゴリ, 係り受け先カテゴリ] を表す。例えば、 $P_n = \{ [機械, ハードウェア] [が, 助詞] \}$  ,  $P_k = \{ [壊れる, 問題] [NULL, NULL] \}$  (n と k には係り受け関係があるものとする) であれば、上のルールを用いて、[ハードウェア, 問題] → [機械, 壊れる] が抽出され、ラベル・概念付きデータベース 180 に登録される。ステップ 470 の判断が Yes であれば処理はステップ 480 で終了する。(以上が係り受け抽出装置 360 の処理である。)



## 【0022】

## [検索・特徴抽出部 130]

図5に検索・特徴抽出部130のブロック図を示す。検索・特徴抽出部130は入力（命令）装置570、表示部510、キーワード検索装置540、命令解析装置520、カテゴリー別キーワード頻度計算装置550、相対頻度計算装置530の各ブロックから構成される。なおキーワード検索装置540及びカテゴリー別キーワード頻度計算装置550は、ラベル・概念付きデータベース180にアクセスして概念情報の検索を行う。好適にはラベル・概念付きデータベース180は、ラベル・概念付きデータに対してインデックスを生成しておき、高速な検索を可能にしておく。

命令解析装置520は、入力（命令）装置570から受け取る命令を解析して、各装置にキーワード・パラメータを送る。入力（命令）装置570は、図14におけるキーボード6、マウス7などに相当し、ユーザの指示に従い、データ分析システムに対して所望の検索、表示を行わせるために使用される。相対頻度計算装置530は、文書の全体、または部分集合に対して相対頻度を計算する装置である。ここで相対頻度とは、全体または任意の集合Xに対して含まれる各キーワードと、任意の集合Yに含まれるキーワードの集合を比較することにより計算される。

## 【0023】

キーワード検索装置540は、命令解析装置520の出力であるキーワードまたはキーワードの組を入力として、文書の全体、または部分集合に含まれるキーワードの数とキーワードを含む文書のIDを得る装置である。該装置によりキーワードを含む文書集合を絞り込むことが出来る。

## 【0024】

カテゴリー別キーワード頻度計算装置550は、命令解析装置520の出力に従い、文書全体、または部分集合に含まれるキーワードの数をカテゴリー別に、頻度順に得る装置である。以下に該装置の出力例を示す。（下の例でINPUTではカテゴリーを指定、N1は[ハードウェア]を表すカテゴリであり、OUTPUTでは[キーワード 出現頻度]となっている。

----- カテゴリー別キーワード頻度計算装置の出力例 -----

[INPUT ] CATEGORY	N1
[OUTPUT] ハードディスク	2033
[OUTPUT] モニタ	1432
[OUTPUT] プリンタ	1001
[OUTPUT] モデム	420
[OUTPUT] スキャナ	212
[OUTPUT] イーサネットカード	143
[OUTPUT] マウス	3

----- カテゴリー別キーワード頻度計算装置の出力例 --終わり--

【 0 0 2 5 】

表示部 5 1 0 は図 6 に示す表示領域（１）と図 7 に示す表示領域（２）を含む G U I 画面から構成される。ユーザが表示部 5 1 0 中に表示された種々の項目を入力（命令）装置 5 7 0 により適宜選択したり、検索のためのパラメータなどを入力することにより、種々の結果（頻度表示、検索結果表示など）を表示部 5 1 0 に表示する。例えば図 6 では、概念 A を横軸、概念 B を縦軸にとった 2 次元の表である。A の列に対して特徴的な B を表すセルが他のセルとは異なる属性で表示されている（例えば概念 A3 の中で特徴的な概念は B1 である。また複数の概念が特徴的なこともある）。上記異なる属性で表示されたセルをクリックすることで、概念 Ax と概念 By を含む集合の検索が行われ、そこで得られた集合に対しても再び検索を行うことができる。またはこの特徴概念の表示を行うことができる。表示部 5 1 0 の図 7 に示す表示領域（２）では、概念 [ソフトウェア] に含まれるキーワードのリストを表示してある。1 つのキーワードに対して 2 つのグラフがあるが、上は頻度、下は相対頻度を表している。さらに頻度・相対頻度のどちらかで並び替えることができる。これらの表示部 5 1 0 における特徴的な概念の抽出及び表示の流れを以下により詳細に説明する。

【 0 0 2 6 】

図 8 に表示領域（１）における、特徴的な概念の抽出及び表示のフローチャートを示す。ステップ 8 2 0 で概念 A と概念 B を選択する。ここで概念 A, B はそれぞれ

れ表示領域（１）でのx, y軸になる。後の計算で、概念Aの要素Axに関して特徴的なものを表示するので、比較したい概念をBにセットする。ステップ840で概念Aを装置540に入力し、概念A中に含まれるキーワードが頻度順に取り出される。また概念Bを装置540に入力し、概念B中に含まれるキーワードが頻度順に取り出される。次にステップ850で、AとBの組合せが終了したかどうかを判断される。終了していなければステップ855で、 $(A_i \ \& \ B_j)$ を装置540に入力し、結果を $P_{ij}$ に保存する。次にステップ860で、正規化が終了したかどうかを判断する。正規化が終了していなければステップ865において、装置530で $B_{Ni}$ を1として正規化する。すなわち  $(P'_{ij} = P_{ij} / B_{Ni})$ となる。図6における表示例中の%表示がそれにあたる。ただし、1つの文章に複数の種類のキーワードが表れるので、これらの値を全部足しても1にはならない。次にステップ870で、相対頻度の計算が終了したかどうかを判断する。計算が終了していなければステップ875で、 $P'_{ij}$ をy軸方向に比較して相対頻度を装置530で調べる。好適には $P'_{ij} (0 < i < a)$ の分布の分散を調べ、各 $P'_{ij}$ が分散の何倍になっているかを相対頻度とする。次に処理はステップ880へ進み、相対頻度の値によって、表示領域（１）の表示色を変え、2次元に表示する。次にステップ890で、目立っている概念A中のキーワードfaと概念B中のキーワードfbの交点をクリックする。最後にステップ898で装置540により、概念Aのfaと概念Bのfbでデータ集合が絞り込まれる。

#### 【0027】

図9に表示領域（２）における、本発明の別の特徴的な概念の抽出及び表示のフローチャートを示す。まずステップ915で表示領域（２）の左側で概念Aを選択する。次にステップ920で装置550により、概念A中に含まれるキーワードが頻度順に取り出される。そしてステップ930で正規化が終了したかどうかを判断される。終了していなければステップ940で相対頻度を計算する。ここでの相対頻度の計算は、あるキーワードが全体または任意の部分集合Xの中に含まれる割合と、任意の部分集合Yに含まれる割合  $Y/X$ で求める。この値が大きいほど、部分集合Yに特徴的なキーワードと考える。ステップ930において正規化が終了した場合ステップ950で、表示領域（２）の右側に頻度と相対頻度を同

時に表示し、頻度または相対頻度で並び替えて表示する。次にステップ 960 で表示領域 (2) の右側でキーワード fa を選択する。そして最後にステップ 970 で入力装置 570 を用いて選択 (クリック) することで、装置 540 による絞り込みが行える。

#### 【0028】

上記 2 つの検索・抽出方法を組み合わせることにより、特徴的な概念を効果的に見つけることが可能となる。例えば表示領域 (2) において、概念 [月] について検索し、「11月」を選択する。次に表示領域 (1) において、[製品名 (コンピュータの機種名)] を縦軸 (比較したい対象)、[問題] を横軸にとる (図 10 参照)。すると「遅い」という [問題] に関して、特徴的な [製品名] が 2 つマークされている。そして、より相対頻度の高い「製品 A」に注目し、「製品 A」と「遅い」の交わっているところを選択 (クリック) する。(図 10 参照) 絞り込まれた状態で表示領域 (2) において、概念 [ハードウェア] について見てみる。すると、頻度の上位 2 番目には「ハードディスク」があり、相対頻度も高く (7.18 倍)、この製品 A に関しては「ハードディスク」が特有の問題を持っていることを推測できる。(図 11 参照)

#### 【0029】

図 13 に本発明で最も特徴的な GUI の例を示す。図 13 において表示領域 (2) の左側で概念 [月] を選ぶ。すると装置 550 により、概念 [月] 中に含まれるキーワードが頻度順に取り出される。そして装置 530 により得られたキーワードの相対頻度が計算され、表示領域 (2) の右側に [月] に含まれる概念が表示される。次に、表示領域 (2) の右側でキーワード「11月」を選択する。これにより装置 540 により、概念 [月] の「11月」でデータ集合が絞り込まれる。

#### 【0030】

表示領域 (1) で概念 [問題] X 軸に設定し、Y 軸に [機種名] をとる。すると装置 550 で概念 [問題] 中に含まれるキーワードが頻度順に取り出される。また装置 550 で概念 [機種名] 中に含まれるキーワードが頻度順に取り出される。これらは検索・正規化され 2 次元に表示される。そして目立っている概念 [問題] の「遅い」と概念 [機種名] の「製品 A」の交点をクリックする。すると装置 540 により

、概念[問題]の「遅い」と概念[機種名]の「製品A」でデータ集合が絞り込まれる。

#### 【0031】

表示領域(2)の左側で概念[ハードウェア]を選ぶ。装置550により、概念[ハードウェア]中に含まれるキーワードが頻度順に取り出される。そして装置530により得られたキーワードの相対頻度が計算される。表示領域(2)の右側に[ハードウェア]に含まれる概念が表示される。結局、上位2番目には「ハードディスク」があり、相対頻度も高く、この製品に関しては、「ハードディスク」が特有の問題ではないかと考えられる。なお上位1番目は、本製品の製品番号などが該当するのでこれらは容易に無視できる。

#### 【0032】

図14に本発明において使用されるデータ分析システムのハードウェア構成例を示す。システム100は、中央処理装置(CPU)1とメモリ4とを含んでいる。CPU1とメモリ4は、バス2を介して、補助記憶装置としてのハードディスク装置13(またはCD-ROM26、DVD32等の記憶媒体駆動装置)とIDEコントローラ25を介して接続してある。同様にCPU1とメモリ4は、バス2を介して、補助記憶装置としてのハードディスク装置30(またはMO28、CD-ROM29、DVD31等の記憶媒体駆動装置)とSCSIコントローラ27を介して接続してある。フロッピーディスク装置20はフロッピーディスクコントローラ19を介してバス2へ接続されている。好適にはラベル付きデータ160、辞書340、ラベル・概念付きデータ370はこれら補助記憶装置に記憶される。

#### 【0033】

フロッピーディスク装置20には、フロッピーディスクが挿入され、このフロッピーディスク等やハードディスク装置13(またはCD-ROM26、DVD32等の記憶媒体)、ROM14には、オペレーティングシステムと協働してCPU等に命令を与え、本発明を実施するためのコンピュータプログラム、オペレーティングシステムのコード若しくはデータを記録することができ、メモリ4にロードされることによって実行される。これらコンピュータ・プログラムのコー

ドは圧縮し、または、複数に分割して、複数の記録媒体に記録することもできる。

#### 【0034】

システム100は更に、ユーザ・インターフェース・ハードウェアを備え、入力をするためのポインティング・デバイス（マウス、ジョイスティック等）7またはキーボード6や、ディスプレイ12を有することができる。好適にはポインティング・デバイス7を用いて、ディスプレイ12に表示された表示領域（1）、表示領域（2）の項目の選択、変更や、パラメータの入力をGUIで行う。また、パラレルポート16を介してプリンタを接続することや、シリアルポート15を介してモデムを接続することが可能である。このシステム100は、シリアルポート15およびモデムまたは通信アダプタ18（イーサネットやトークンリング・カード）等を介してネットワークに接続し、他のコンピュータ、サーバ等と通信を行う。本発明にデータ分析システムは、必要なデータベースを、通信回線を介して接続された外部のサーバ、WAN、LAN内にあるローカルなサーバなどに記憶してもよい。どちらにしても本発明の実施の制限となるものではない。またシリアルポート15若しくはパラレルポート16に、遠隔送受信機器を接続して、赤外線若しくは電波によりデータの送受信を行ってもよい。

#### 【0035】

スピーカ23は、オーディオ・コントローラ21によってD/A（デジタル／アナログ変換）変換されたサウンド、音声信号を、アンプ22を介して受領し、サウンド、音声として出力する。また、オーディオ・コントローラ21は、マイクロフォン24から受領した音声情報をA/D（アナログ／デジタル）変換し、システム外部の音声情報をシステムにとり込むことを可能にしている。ViaVoice（IBM商標）などのアプリケーションを用いて、本発明のGUIコマンド部の操作を音声コマンドによる操作で代用してもよい。さらにホームページ・リーダー（IBM商標）などアプリケーションを用いて、表示された特徴的概念を有する検索結果などを音声で読み上げるようにしてもよい。

#### 【0036】

このように、本発明のデータ分析システムは、通常のパーソナルコンピュータ

(PC) やワークステーション、ノートブックPC、パームトップPC、ネットワークコンピュータ、コンピュータを内蔵したテレビ等の各種家電製品、通信機能を有するゲーム機、電話、FAX、携帯電話、PHS、電子手帳、等を含む通信機能有する通信端末、または、これらの組合せによって実施可能であることを容易に理解できるであろう。ただし、これらの構成要素は例示であり、その全ての構成要素が本発明の必須の構成要素となるわけではない。

【0037】

【発明の効果】

本発明により、大量データから有効な知識を獲得する方法およびシステムが提供される。また有効な知識を獲得するための、ユーザビリティに優れたGUIによる分析方法およびシステムが提供される。

【図面の簡単な説明】

【図1】

本発明のデータ分析システム概要を示すブロック図である

【図2】

データ変換部のフローチャートである。

【図3】

概念抽出部のブロック図である。

【図4】

概念抽出部のフローチャートである。

【図5】

検索・特徴抽出部のブロック図である。

【図6】

表示部における表示領域(1)の例である。

【図7】

表示部における表示領域(2)の例である。

【図8】

表示領域(1)における概念の抽出及び表示のフローチャートである。

【図9】

表示領域（２）における概念の抽出及び表示のフローチャートである。

【図 1 0】

表示領域（１）の典型的な表示例である。

【図 1 1】

表示領域（２）の典型的な表示例である。

【図 1 2】

概念抽出部の処理を具体的な文を用いて説明した図である

【図 1 3】

表示領域（１）と表示領域（２）を含む G U I 画面の表示例である。

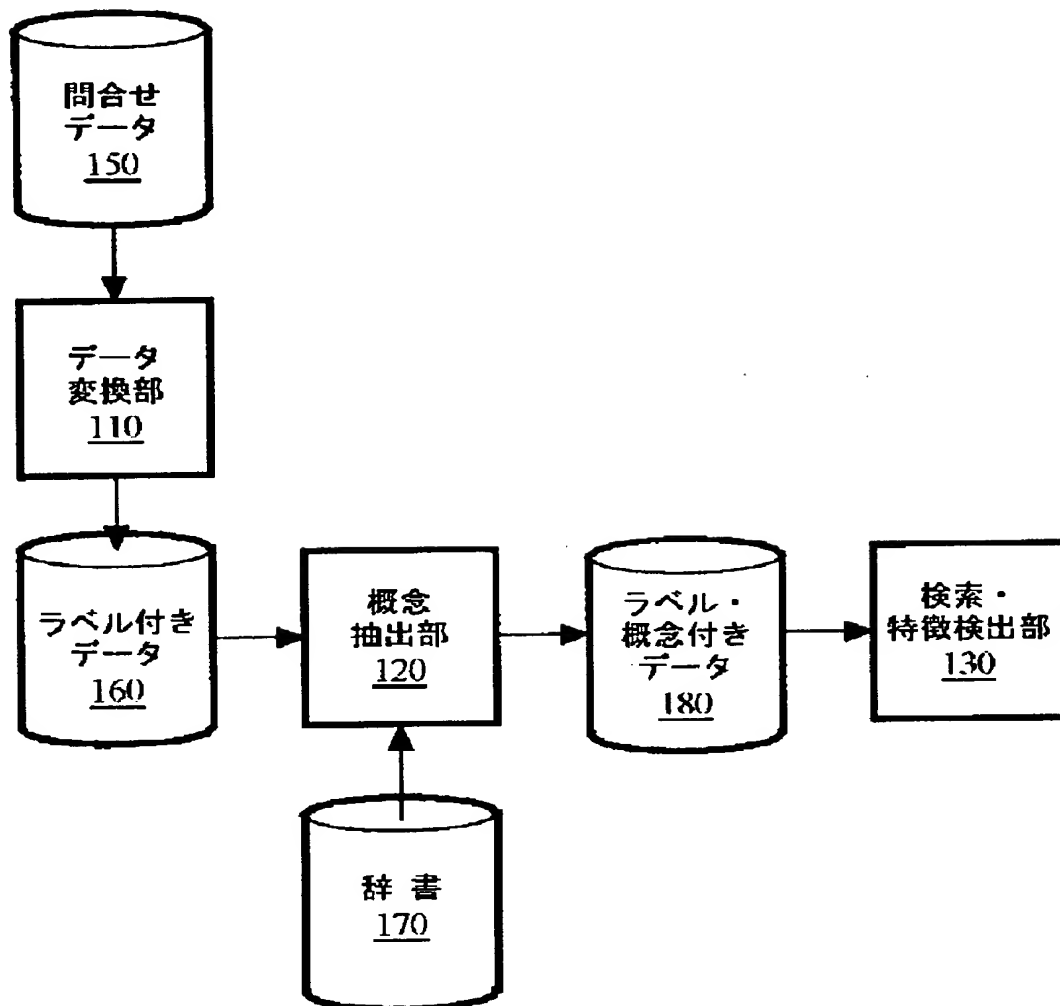
【図 1 4】

本発明で用いるハードウェア構成の実施例である。

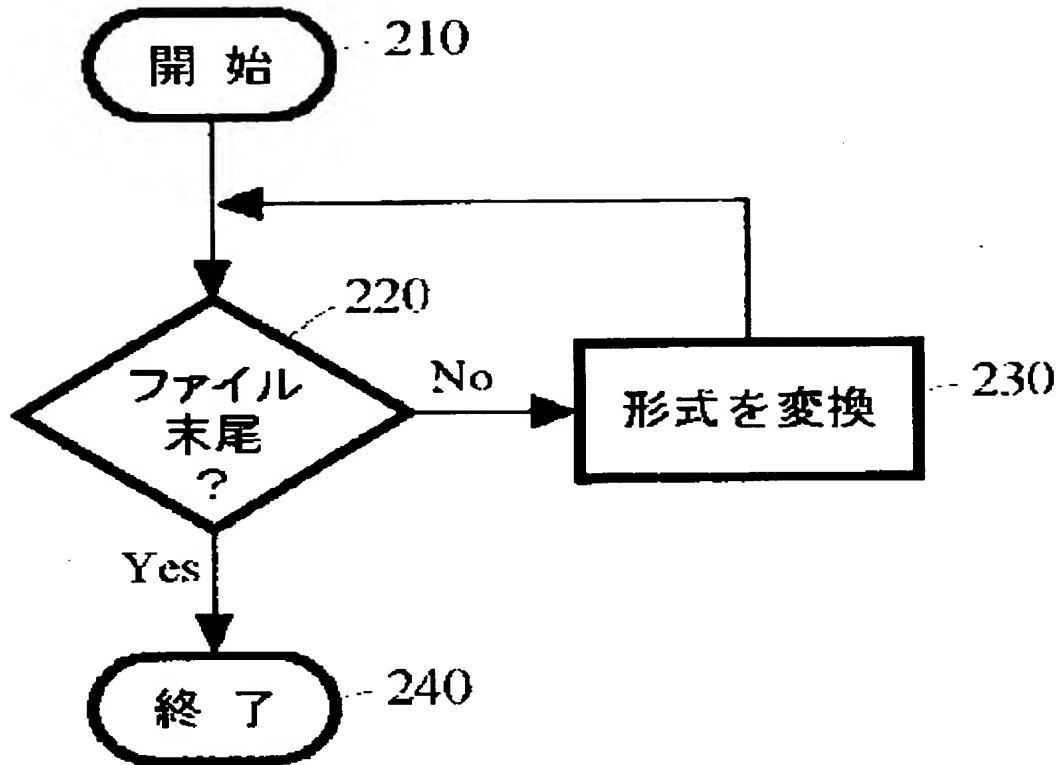


【書類名】 図面

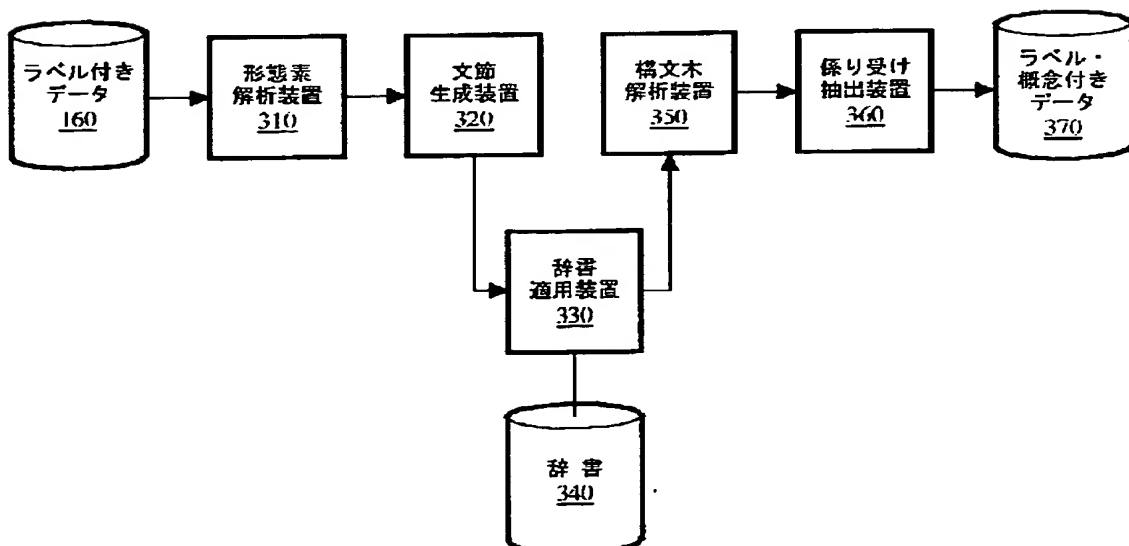
【図 1】



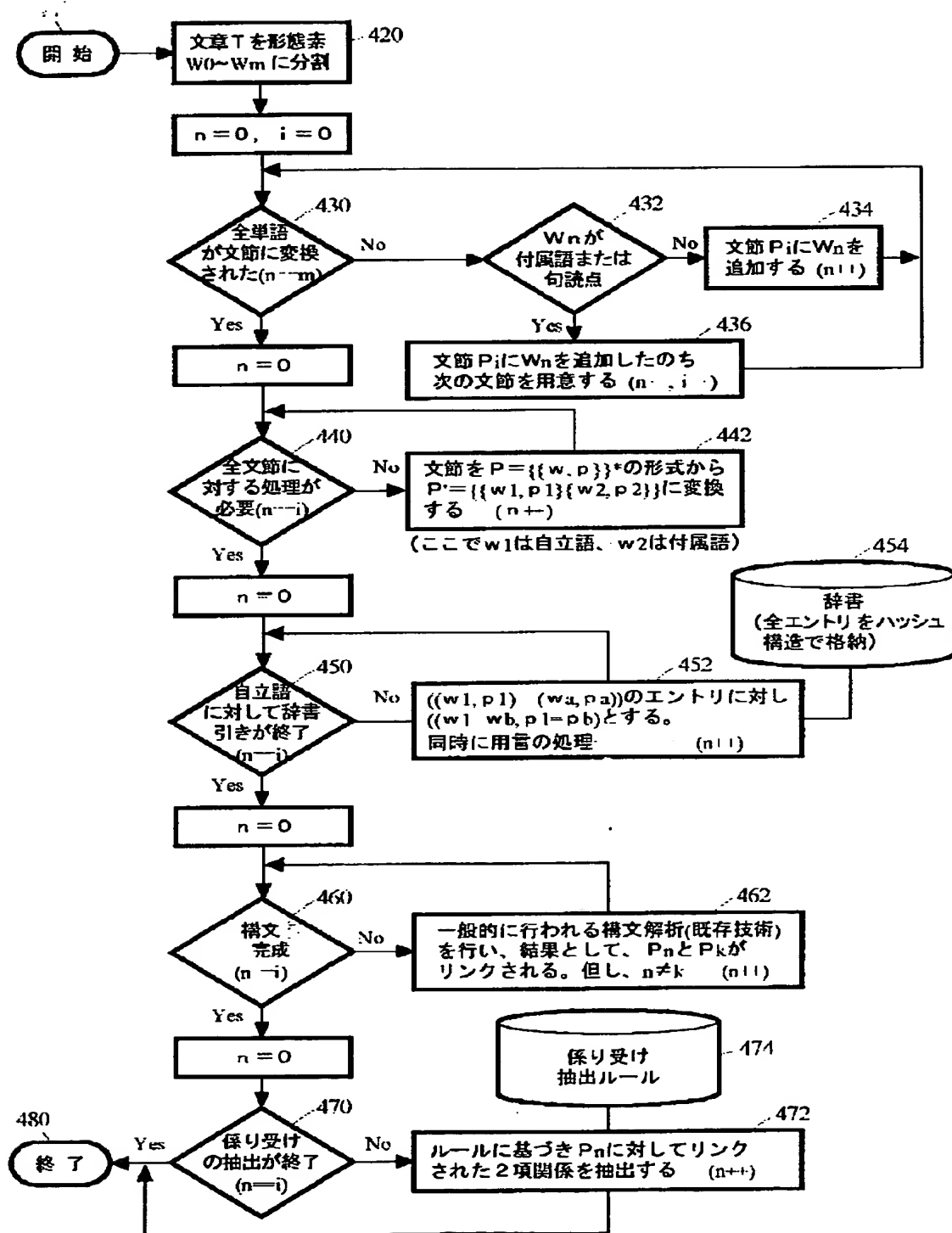
【図 2】



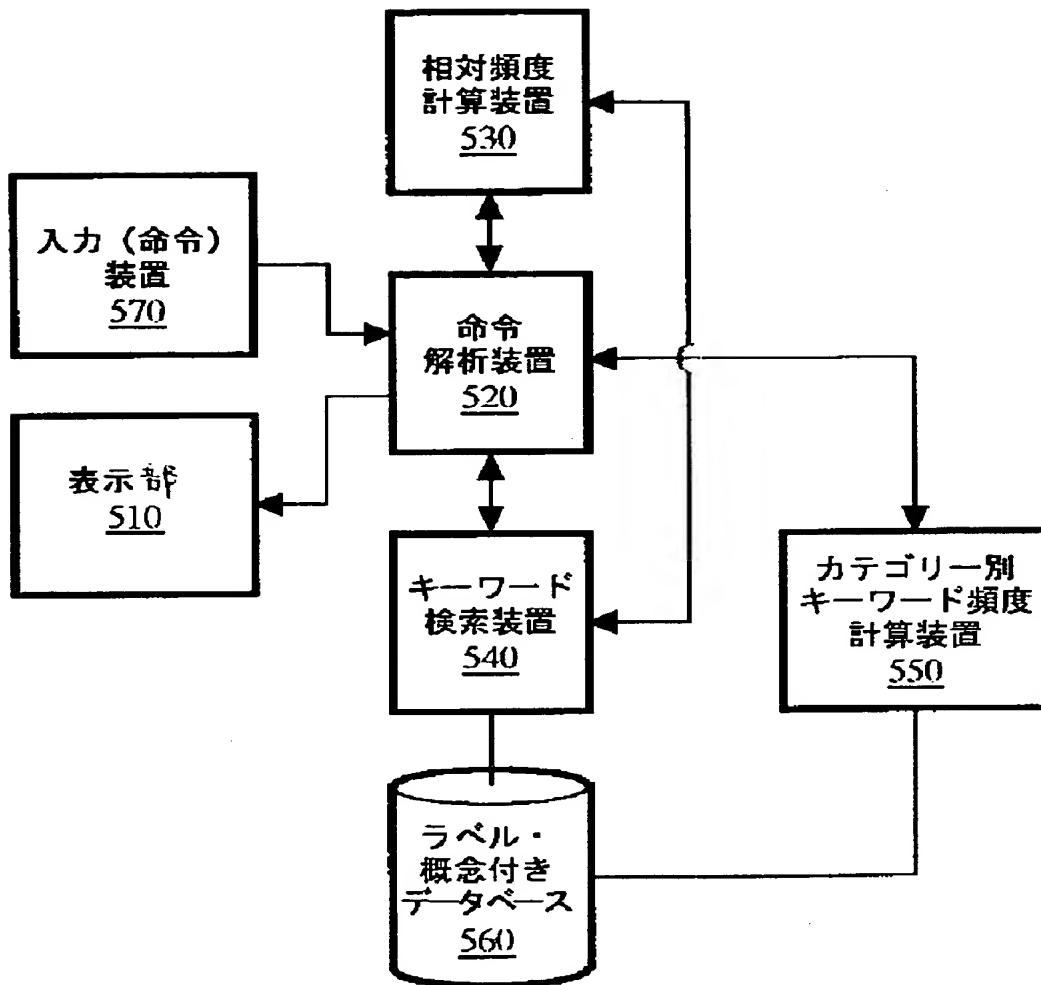
【図 3】



【図 4】



【図 5】



【図 6】

		概念 A						
		A1	A2	A3	A4	A5	A6	A7
概念 B	B1	1720 (89.6%)	56 (2.9%)	118 (6.1%)	7 (0.3%)	8 (0.4%)	7 (0.3%)	1 (0.0%)
	B2	1255 (90.5%)	21 (1.5%)	100 (7.2%)	4 (0.2%)	2 (0.1%)	2 (0.1%)	2 (0.1%)
	B3	1072 (95.0%)	28 (2.4%)	18 (1.5%)	5 (0.4%)	0 (0.0%)	3 (0.2%)	1 (0.0%)
	B4	501 (90.1%)	9 (1.6%)	40 (7.1%)	1 (0.1%)	3 (0.5%)	2 (0.3%)	0 (0.0%)
	B5	649 (93.2%)	26 (3.7%)	17 (2.4%)	1 (0.1%)	0 (0.0%)	3 (0.4%)	0 (0.0%)
	B6	682 (95.5%)	13 (1.8%)	15 (2.1%)	3 (0.4%)	0 (0.0%)	1 (0.1%)	0 (0.0%)
	B7	676 (96.9%)	11 (1.5%)	5 (0.7%)	3 (0.4%)	1 (0.1%)	1 (0.1%)	0 (0.0%)
	B8	673 (90.7%)	34 (4.5%)	26 (3.5%)	4 (0.5%)	0 (0.0%)	5 (0.6%)	0 (0.0%)
	B9	528 (74.3%)	134 (18.8%)	10 (1.4%)	14 (1.9%)	2 (0.2%)	22 (3.0%)	0 (0.0%)
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	

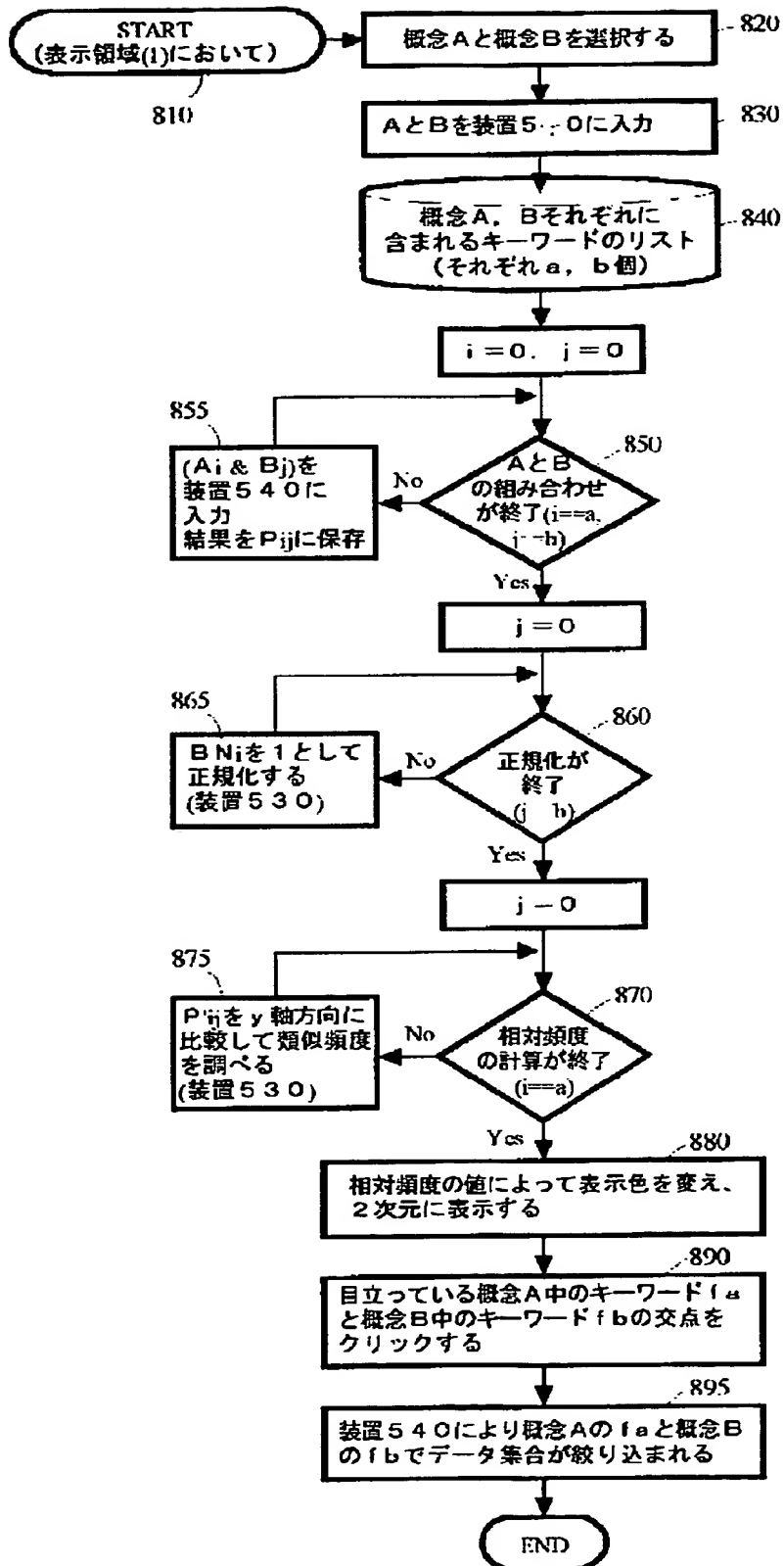
表示領域 ( 1 )

【図 7】

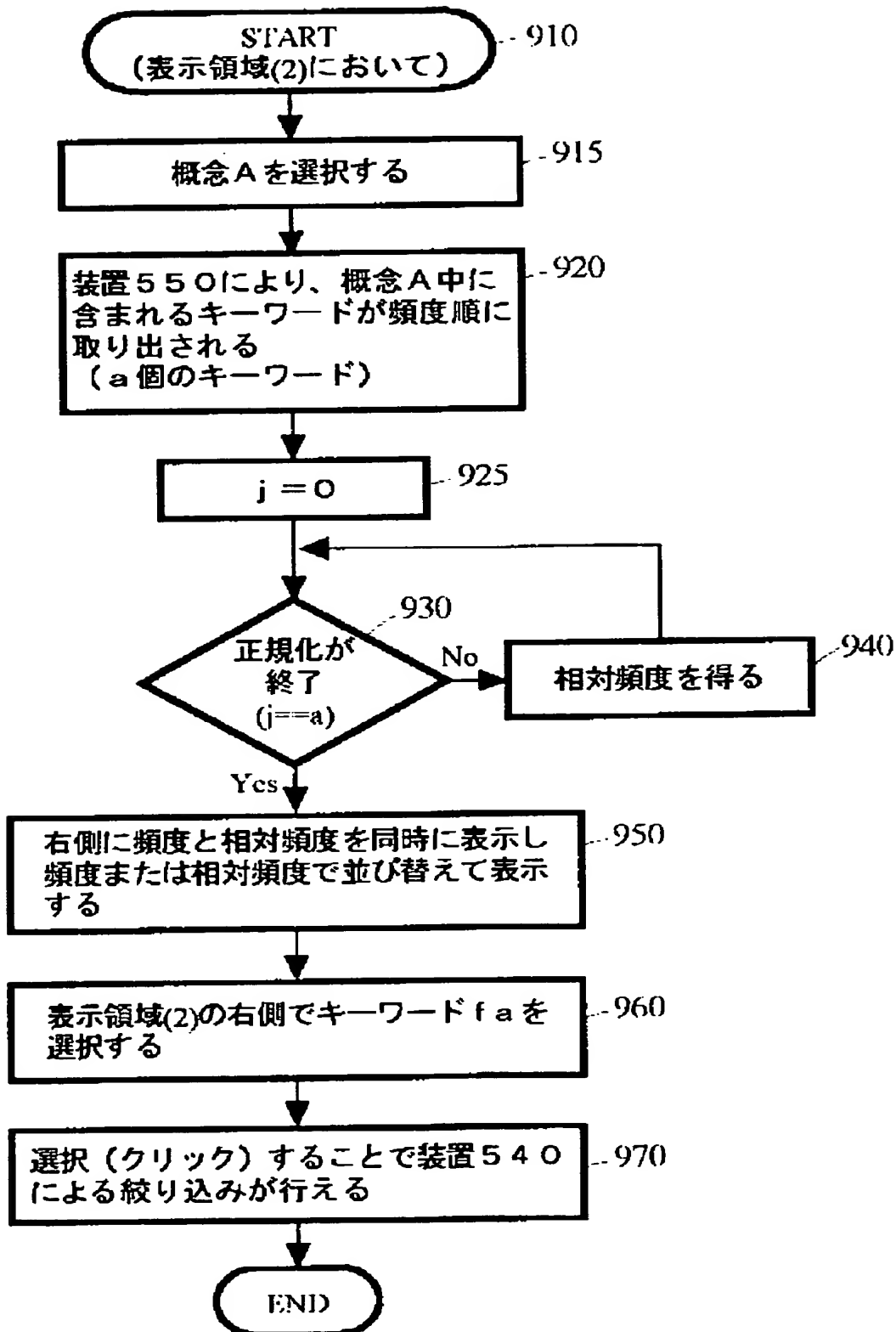
件数・相対頻度		ソート ○ 件数 ◎ 相対		スナップショット
概念 B	ソフトウェア	120	0.0	WINDOWS98
	ハードウェア	118	0.0	WINDOWS98
	専門用語	77	0.0	アップグレード
	コマンド	53	0.0	リカバリー
	対象コンポ	43	0.0	メッセージ
	問題種別	34	0.0	WINDOWS
	CALL 種別	27	0.0	ドライバー
	回答・対応種別	27	0.0	ソフト
	機種名	27	0.0	起動時
	解決時間 (分)	19	0.0	アプリケーション
	CALL 回数			
	対応担当者数			
	対応チーム数			
	名詞			
	固有名詞			
	その他			
	人名			
	組織名			

表示領域 ( 2 )

【図 8】



【図9】



【図 1 0】

						遅い	
	19 (2.58%)	6 (0.82%)	9 (1.22%)	6 (0.82%)	2 (0.27%)	2 (0.27%)	3 (0.41%)
	18 (2.02%)	7 (0.79%)	18 (2.02%)	4 (0.45%)	4 (0.45%)	3 (0.34%)	3 (0.34%)
	21 (2.75%)	14 (1.83%)	15 (1.97%)	12 (1.57%)	12 (1.57%)	14 (1.83%)	7 (0.92%)
	57 (5.76%)	12 (1.21%)	11 (1.11%)	8 (0.81%)	3 (0.31%)	1 (0.10%)	5 (0.51%)
	55 (4.8%)	20 (1.75%)	7 (0.61%)	11 (0.96%)	2 (0.17%)	1 (0.09%)	5 (0.44%)
製品 A	20 (2.58%)	15 (1.94%)	14 (1.81%)	8 (1.03%)	6 (0.78%)	20 (2.58%)	6 (0.78%)
	15 (2.11%)	18 (2.53%)	12 (1.69%)	9 (1.27%)	1 (0.14%)	5 (0.70%)	3 (0.42%)
	17 (2.29%)	10 (1.35%)	33 (4.45%)	6 (0.81%)	1 (0.13%)	2 (0.27%)	2 (0.27%)
	15 (2.76%)	7 (1.29%)	8 (1.47%)	5 (0.92%)	2 (0.37%)	0 (0.00%)	3 (0.55%)
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

表示領域 ( 1 )

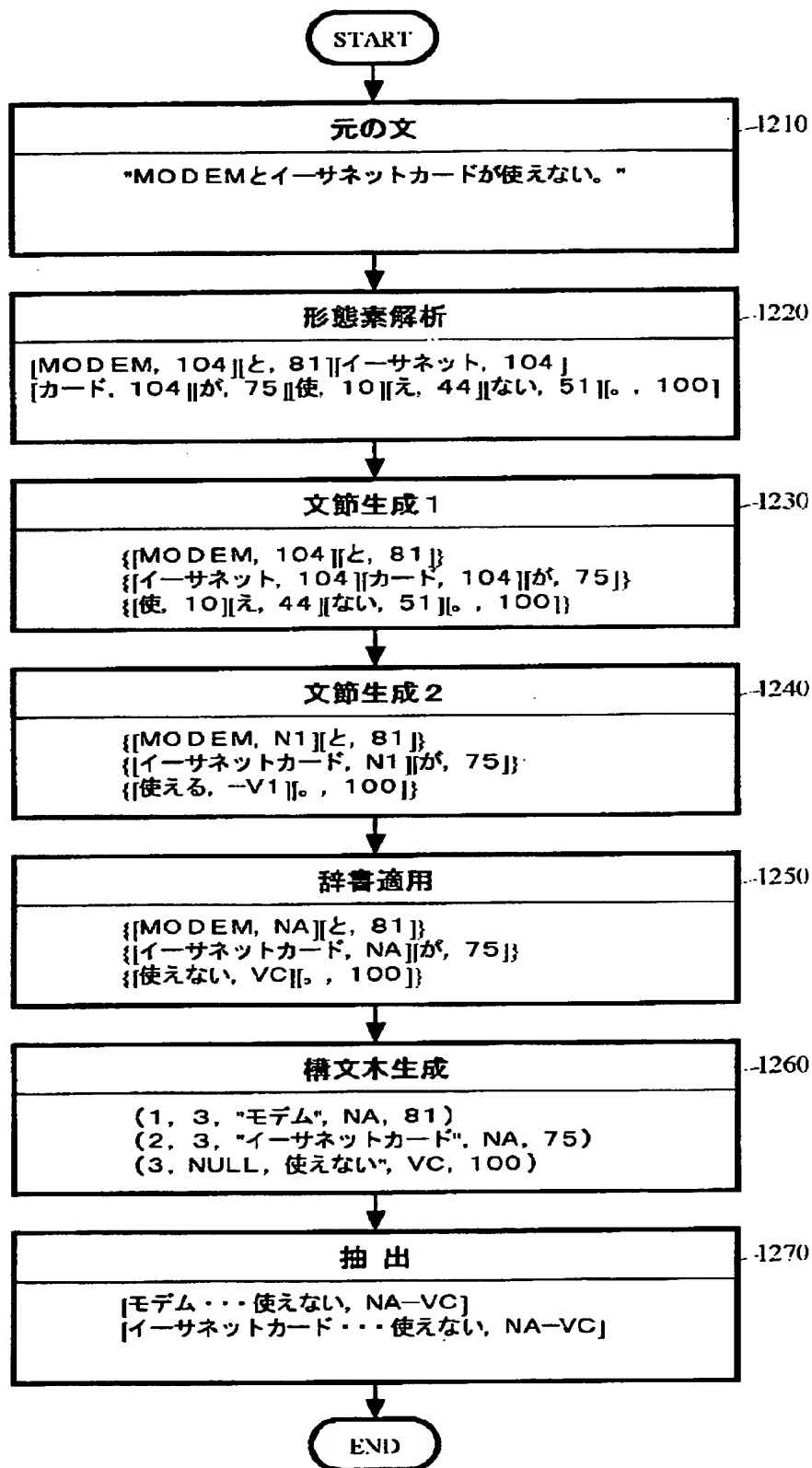
【図 1 1】

件数・相対頻度	ソート	● 件数 ○ 相対 ○ 項目	スナップショット
ソフトウェア	1.37	-----	△
ハードウェア	1.18	ハードディスク	
専門用語	3.1	-----	
コマンド	2.0	-----	
対象コンポ	2.0	-----	
問題種別	2.0	-----	
CALL 種別	1.0	-----	
回答・対応種別	1.0	-----	
機種名	1.0	-----	
解決時間 (分)	1.0	-----	
CALL 回数	1.0	-----	
対応担当者数	1.0	-----	
対応チーム数	1.0	-----	
名詞	1.0	-----	
固有名詞	1.0	-----	
その他		-----	▽
人名			▽
組織名			▽

表示領域 ( 2 )



【図 1 2】



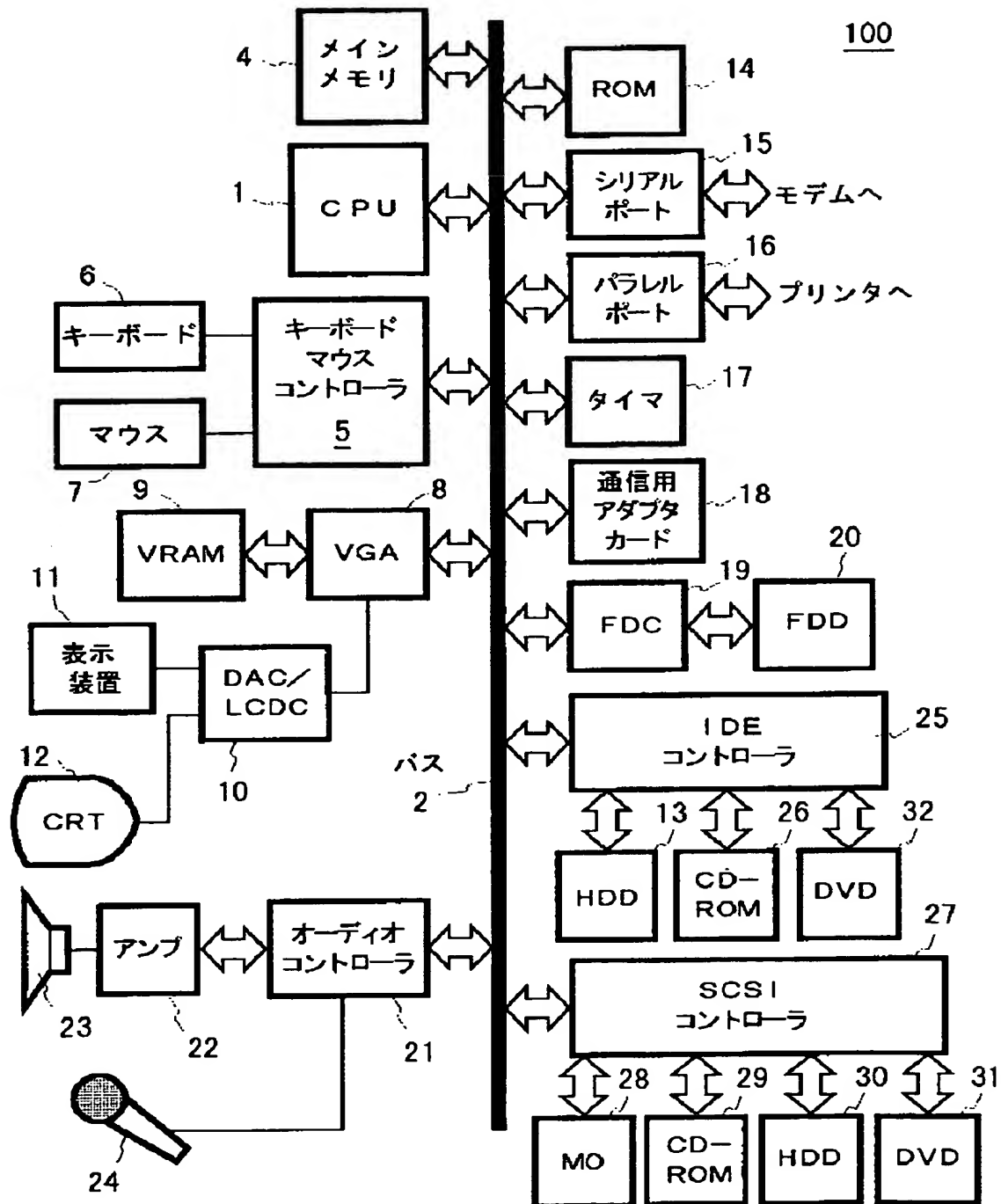
【図 13】

サンプリング件数 1000 データベース名		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
キーワード入力		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
履歴名		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
削除 クリア	履歴を開く 履歴の保存	件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
36758 MONTH(1998/08)		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
841 遅延		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
29 遅い+時間		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
スナップショット		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
19980801 起動に異常に時間がかかる		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
19980801 起動にもものすごく時間がかかる		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
19980802 起動に30分かかる		件数・相対頻度		ソート ● 件数 ○ 相対 ○ 項目		スナップショット	
2Dグラフ		増減グラフ		トピック抽出 時系列グラフ		スナップショット	
ソフトウェア	ソフトウェア	起動しない	止まる	起動できない	できない	おかしい	不明
ハードウェア	ハードウェア	製品A 37(4.3%)	21(2.4%)	13(2.1%)	14(1.6%)	12(1.4%)	4(0.4%)
専門用語	専門用語	製品B 36(5.1%)	20(2.8%)	21(2.9%)	13(1.8%)	7(0.9%)	6(0.6%)
コマンド	コマンド	製品C 22(3.2%)	16(2.3%)	18(2.8%)	10(1.4%)	11(1.1%)	1(0.1%)
対象コンボ	対象コンボ	製品D 12(2.3%)	10(1.9%)	20(3.8%)	3(0.5%)	9(1.7%)	2(0.3%)
問題種別	問題種別	製品E 33(4.9%)	27(4.0%)	13(1.9%)	8(1.2%)	17(2.5%)	4(0.6%)
CALL種別	CALL種別	製品F 20(3.5%)	28(4.9%)	10(1.7%)	7(1.2%)	14(2.4%)	1(0.1%)
回答・対応種別	回答・対応種別						
機種名	機種名						
解決時間(分)	解決時間(分)						
CALL回数	CALL回数						
対応担当者数	対応担当者数						
対応チーム数	対応チーム数						
名詞	名詞						
固有名詞	固有名詞						
その他	その他						
人名	人名						
組織名	組織名						

表示領域  
(1)

表示領域  
(2)

【図 14】



【書類名】 要約書

【要約】

【課題】 大量データから有効な知識を獲得する方法およびシステムを提供することである。

【解決手段】

本発明は文書を含む大量のデータから特異な特徴を有する概念を抽出することにより、有効な知識を獲得する方法およびシステムを提供する。本発明は概念抽出装置と特徴的概念抽出装置と有する。前記概念抽出は文書データを含むデータからカテゴリ別の概念を抽出する。前記特徴的概念抽出装置は前記抽出した概念の中から特徴的な概念を抽出する装置であって、前記カテゴリ別の概念において、同一カテゴリに属する概念のうち、対応する別のカテゴリに属する概念の中で占める割合が既定値を超えている概念を抽出する。

【選択図】 図 1

認定・付加情報

特許出願の番号	平成11年 特許願 第192618号
受付番号	59900650449
書類名	特許願
担当官	第七担当上席 0096
作成日	平成11年 7月 9日

<認定情報・付加情報>

【提出日】	平成11年 7月 7日
-------	-------------

出 願 人 履 歴 情 報

識別番号 [390009531]

1. 変更年月日 1990年10月24日  
[変更理由] 新規登録  
住 所 アメリカ合衆国10504、ニューヨーク州 アーモンク (番地なし)  
氏 名 インターナショナル・ビジネス・マシーンズ・コーポレイション